051

052

053

054

055

056

057

# GraphPad: Inference-Time 3D Scene Graph Updates for Embodied Question Answering

Anonymous CVPR submission

# Paper ID 8

#### Abstract

Structured scene representations are a core component 001 of embodied agents, helping to consolidate raw sensory 002 streams into interpretable, modular, and searchable for-003 mats. Due to their high computational overhead, many ap-004 proaches build such representations in advance of the task. 005 006 However, when the task specifications change, such static approaches become inadequate as they may miss key ob-007 jects, spatial relations, and details. We introduce Graph-008 **Pad**, a modifiable structured memory that an agent can tai-009 lor to the needs of the task through API calls. It comprises a 010 011 mutable scene graph representing the environment, a navigation log indexing frame-by-frame content, and a scratch-012 pad for task-specific notes. Together, GraphPad serves as 013 a dynamic workspace that remains complete, current, and 014 aligned with the agent's immediate understanding of the 015 scene and its task. On the OpenEQA benchmark, Graph-016 017 Pad attains 55.3% accuracy-+3.0 pp over an image-only baseline using the same vision-language model-while op-018 erating with five times fewer input frames. These results 019 show that allowing online, language-driven refinement of 020 021 3-D memory yields more informative representations with-022 out extra training or data collection.

# **023 1. Introduction**

A household robot has just scanned the kitchen when a user 024 025 asks, "Is the red mug back inside the upper cupboard?" Because the scene graph was built from only a handful of 026 earlier keyframes, it lacks both a node for the mug and the 027 relation linking it to the cupboard interior. With this criti-028 029 cal information missing, the agent must either guess, rescan the entire scene, or refuse the query altogether. This sce-030 nario highlights a fundamental limitation in current embod-031 ied AI: structured 3D memories are typically finalized 032 before the downstream task is known, frequently omit-033 ting objects and spatial relations that later prove essen-034 035 tial for action or reasoning.

Inherent limitations of static memories. Modern em-036 bodied systems commonly pair a vision-language model 037 (VLM) with a pre-computed 3D scene graph or semantic 038 map for symbolic reasoning [11, 12]. These graphs com-039 press dense RGB-D streams into discrete entities and re-040 lations that a VLM can leverage for navigation, manipu-041 lation, or question answering. However, key compression 042 decisions-which keyframes to store, which detections to 043 accept, which relations to represent as edges-are made 044 without knowledge of the eventual user request. Conse-045 quently, graphs often lack precisely the objects or spatial 046 relations required by a subsequent task, a failure mode well-047 documented in embodied question answering research [5, 048 15]. 049

Current solutions follow two primary approaches. *Overprovisioning* retains nearly every candidate detection, incurring substantial memory and computational overhead [2]. *Offline enhancement pipelines* append heuristically chosen affordances or relations for anticipated tasks [12], but each new domain demands additional engineering effort and cannot recover information never initially detected. Neither strategy scales effectively when user requests vary widely.

Our approach: language-guided memory updates. We 058 propose that an embodied agent should be able to dynam-059 ically revise its structured scene memory when reasoning 060 reveals a gap in its knowledge. Rather than discarding an 061 inadequate structured representation-and thereby losing 062 its benefits-or recreating it from scratch at high compu-063 tational cost, the agent should be able to modify and en-064 hance the existing representation. This allows the agent to 065 incorporate task-relevant information without starting over. 066 Moreover, the agent should be able to do this through tar-067 geted perceptual queries on specific keyframes, and seam-068 lessly integrate the newly discovered information into its 069 scene graph. 070

**GraphPad.** We introduce **GraphPad**, a modifiable 3D 071 scene graph memory whose content can be updated at in-

133

134

135

136

137

138

139

140

146

147

148

149

150

151

152

153

154

155

156

157

158

159

160

161

162

163

164

165

166

167

ference time through language-level commands issued by 073 074 the same VLM that performs high-level reasoning. Whenever uncertainty arises or missing information is identi-075 fied, the VLM can: (i) retrieve additional keyframes from 076 077 its navigation log, (ii) insert previously undetected objects or spatial relationships, and (iii) annotate existing nodes 078 with task-specific attributes. These operations are exposed 079 as language-callable functions that the reasoning agent it-080 081 self can invoke, eliminating the need for task-specific reasoning code.<sup>1</sup> Through successive targeted updates, the 082 083 graph evolves from a coarse initial representation into a task-conditioned model containing precisely the informa-084 tion needed for the current query or action planning. 085

#### **086** Contributions.

- We formulate *language-driven online editing* of structured 3D memories as a solution to the task-memory mismatch inherent in fixed scene graphs.
- We present *GraphPad*, a system in which a single VLM
   both identifies knowledge gaps and updates its 3D repre sentation via language function calls during inference.
- On the OpenEQA benchmark [15], we demonstrate that GraphPad improves spatial question answering from 52.3% to 55.3% while reducing the number of processed frames from 25 to 5, outperforming both image-only and static-graph baselines without additional training.

098By recasting perception as an interactive, language-099mediated process rather than a static preprocessing step,100GraphPad enables more efficient and effective reasoning101about complex 3D environments—a critical capability for102agents that must bridge language understanding with phys-103ical action.

#### **104 2. Background and Related Work**

Current embodied AI systems struggle with a fundamental
limitation: their memory representations often lack critical
information needed during task execution. This section examines key research areas relevant to this challenge.

# **109 2.1. Embodied Memory Representations**

Frame-level memories store raw RGB-D streams with
learned descriptors. Systems like ReMEmbR [2] and
Embodied-RAG [28] index thousands of images, allowing
retrieval of relevant observations during inference. While
flexible, these approaches face computational overhead
from processing large image collections and delegate all
spatial reasoning to language models.

Semantic metric maps address these limitations by embedding visual features within geometric reconstructions.

VL-Maps [10], PLA [6], and OpenScene [16] support spa-<br/>tial grounding but lack explicit object boundaries and rela-<br/>tionships essential for symbolic reasoning.119120

**3D scene graphs** compress perception into discrete122entities and relations. Hydra [11] pioneered real-time123scene graph construction, with extensions for monocular124inputs [1], multi-robot fusion [7], and open-vocabulary la-125beling [8]. These structures align with structured VLM126prompting but remain largely static after construction.127

#### 2.2. Memory Update Mechanisms

Several existing methodologies employ updatable memo-<br/>ries; however, they typically fall short by not dynamically<br/>refining scene graphs with new semantic information cru-<br/>cial for the task at hand but absent in the initial recording:129<br/>130131<br/>132

**Dynamic object tracking** systems like DynaMem [13], OpenIN [24], and Kimera [1] update spatial relations when objects move but rely on predefined update policies rather than reasoning-triggered modifications.

**Incremental discovery** methods like Moma-LLM [9] and Search3D [23] can integrate new objects but employ fixed detection heuristics rather than responding to specific reasoning gaps.

Memory selection approaches like 3D-Mem [29] and141KARMA [25] retain informative views but cannot add142structure absent from initial observations.Similarly,expanded-context approaches [3] can reference more images but cannot insert missing relations.143

#### 2.3. Task-Memory Alignment

Ensuring memory representations contain task-relevant information remains a central challenge:

**Graph adaptation** techniques attempt to tailor existing memories to specific tasks. Information-theoretic approaches [14] compress graphs while preserving taskrelevant nodes. SayPlan [20] and SayNav [18] dynamically contract graphs during planning. Neither approach handles adding new information. LLM-enhanced scene graphs [12] augments an existing scene graph with task-specific affordances for the household rearrangement task. However, they rely on a complex scene graph enhancement pipeline that is domain-specific and not adaptable to diverse, unforeseen downstream tasks.

**Image-level reasoning** systems like Bumble [22] and TagMap [30] provide direct visual access but sacrifice explicit relational structure. Graph-based reasoning methods [28] provide structured context but treat graphs as readonly, requiring new exploration rather than memory editing.

Empirical analyses [26] confirm that missing nodes and relations—not reasoning failures—frequently cause performance errors in embodied tasks.

<sup>&</sup>lt;sup>1</sup>Generic vision modules—object detector, mask extractor, depth backprojection—remain necessary but are *not* specialized for any particular downstream task.

**CVPR** 



Figure 1. Overview of GraphPad for embodied question answering. The top row illustrates the Structured Scene Memory (SSM) components: Scene Graph with associated Graphical Scratch Pad, Frame Memory containing sparse keyframes, and Navigation Log indexing frame metadata. The bottom row demonstrates the inference process: given the question about a white object above the TV, the VLM first examines the initial memory state (left), identifies missing information and calls the analyze\_frame API on a promising frame (middle), then integrates the newly detected air conditioner into the scene graph and scratch pad before providing the answer (right). This process enables dynamic, task-specific memory updates without exhaustive preprocessing.

#### 168

#### **2.4. Embodied Question Answering**

Embodied Question Answering (EQA) benchmarks offer a
standardized evaluation of memory adequacy. The episodic
memory variant (EM-EQA) [15] provides agents with fixed
observations, isolating memory representation quality from
exploration.

Recent approaches to EM-EQA highlight representation
completeness as a critical factor. GraphEQA [21] uses semantic graphs for viewpoint selection, while 3D-Mem [29]
curates informative snapshots. Both systems nevertheless
fail when queried entities are absent from their representations.

GraphPad addresses the task-memory alignment prob-180 181 lem through targeted scene graph updates during inference. Unlike previous work, our approach enables the reasoning 182 VLM to detect knowledge gaps and modify its own mem-183 ory representation through three specific operations: frame 184 retrieval, entity/relation insertion, and semantic annotation. 185 This maintains the efficiency advantages of structured rep-186 187 resentations while addressing their typically static nature.

# 3. Methodology

GraphPad builds a *Structured Scene Memory* (SSM) from a sparse set of RGB-D keyframes and gives the visionlanguage model (VLM) that answers questions three callable functions to *edit* that memory during inference. We first describe the **agentic reasoning loop** (Sec. 3.1), then the four data structures that constitute the SSM (Sec. 3.2), and finally the **Modifiability APIs** (Sec. 3.3). 199

Throughout we follow the episodic-memory EQA protocol [15]: the agent is given every k-th RGB-D frame of a pre-recorded scan together with camera poses; no new sensing is possible after deployment.

# 3.1. Agentic Reasoning Loop

At test time, the VLM receives both the initial SSM and a<br/>natural language query q. The reasoning process unfolds201iteratively: the VLM analyzes what information it needs to<br/>answer the question and repeatedly invokes Modifiability203APIs to augment its understanding of the scene.205

In each iteration, the VLM examines the current scene 206 memory, identifies knowledge gaps relevant to the question, 207

188

196

197

198

199

200

275

276

277

278

279

280

281

282

283

284

285

286

287

288

289

290

291

292

293

294

295

296

297

and selects: (1) which API to call, (2) which frame to analyze, and (3) what specific information to seek (expressed as a natural language query). The selected API returns new objects, relations, or semantic annotations that are integrated into the scene memory. This process continues until the VLM determines it has sufficient information or reaches a maximum of *m* allowed API calls.

A critical constraint ensures grounded responses: the VLM must support every factual claim in its final answer with dual evidence—at least one frame from Frame Memory and corresponding notes from the Scratch-Pad. This constraint ensures that the VLM adequately utilizes its scene memory and searches the scene thoroughly prior to responding.

This approach fundamentally differs from static memory systems: instead of relying on pre-built representations to contain all necessary information, GraphPad actively identifies and fills knowledge gaps during reasoning. In our experiments, most questions are answered with just 2-3 targeted API calls, demonstrating efficient reasoning without exhaustive scene analysis.

#### **3.2. Structured Scene Memory**

230 The SSM contains four mutually linked structures:

#### 231 Scene Graph

A directed multigraph  $G = (\mathcal{N}, \mathcal{E})$  whose nodes rep-232 resent object *tracks*. Each node  $n_i$  stores a point cloud 233  $P_i$ , pooled visual embedding  $V_i$ , pooled language em-234 bedding  $L_i$ , caption  $C_i$ , room/floor ID, and a list of 235 keyframes in which the object is visible. Edges encode 236 four view-invariant spatial relations critical for planning 237 238 and manipulation: on top of, subpart of, contained in, and attached to. 239

#### 240 Graphical Scratch-Pad

241 Mirrors  $\mathcal{N}$  but adds a free-form notes field initialized 242 empty; the APIs write task-specific information here 243 during reasoning.

#### 244 Frame Memory

245 An initial set of  $n_{img}$  evenly spaced keyframes. Addi-246 tional frames requested by the APIs are appended (no 247 eviction is used in our experiments).

#### 248 Navigation Log

For each keyframe: room, textual field-of-view tag, egocentric motion label (from pose deltas), and the set of
visible node IDs. The log serves as a structured index,
guiding the VLM in selecting candidate frames likely
to contain information about specific objects or spatial
relationships.

**Initial construction.** For every k-th RGB image  $I_t$ , we 255 run a VLM detector that outputs bounding boxes and cap-256 tions. Each box is passed to SAM to obtain a mask; the 257 mask is back-projected with depth into a point cloud and 258 voxel-downsampled to 0.02 m. Noise is removed by keep-259 ing only the largest DBSCAN cluster (default sklearn 260 parameters). Visual embeddings are extracted using CLIP 261 ViT-L/14; language embeddings come from BGE [27]. 262

**Track association.** A new detection  $D_i$  is matched to an existing track  $T_j$  when the vote 264

$$S_{ij} = \mathbf{1}[V_i \cdot V_j > 0.7] + \mathbf{1}[L_i \cdot L_j > 0.8] + \mathbf{1}[G_{ij} > 0.4]$$
(1) 265

exceeds 2, where  $G_{ij}$  is the fraction of points in  $D_i$  within 266  $\delta_g$  of  $T_j$  (we use  $\delta_g = 5$  cm). This voting scheme requires 267 at least two of the following conditions to hold with suffi-268 cient confidence: visual feature similarity, caption semantic 269 similarity, or spatial overlap. Visual and language embed-270 dings of the matched track are updated by an exponential 271 moving average with  $\alpha = 0.5$ ; unmatched detections start 272 new tracks. 273

Edge discovery. Every three frames we prompt the VLM with the current frame plus the JSON list of visible  $\langle bbox, caption \rangle$  pairs; the model predicts all pairwise relations among the four relation types. Each predicted edge is stored with a subject-ID, object-ID, relation label, and the VLM's free-form justification string.

**Caption consolidation.** Accumulated captions on a track are periodically compressed by prompting the VLM with the list and asking for a single sentence that faithfully paraphrases all entries.

**Room/floor labels.** We adopt the HOV-SG pipeline [26]: floors from height-histogram modes and rooms via water-shed segmentation on wall skeletons; room labels derive from CLIP similarity to a fixed set of class names.

When the SSM is passed to the VLM for reasoning, the Scene Graph, Scratch-Pad, and Navigation Log are serialized to JSON, while Frame Memory is supplied as interleaved images with their Frame IDs.

# **3.3. Modifiability APIs**

All three APIs receive a *Frame ID* and a natural-language *query*. Each returns a JSON patch containing new nodes, edges, or scratch-pad notes plus evidence pointers that are incorporated into the SSM.

#### find\_objects

Detects previously unseen instances relevant to *query* in the specified frame and fuses them into G. The function leverages the VLM's bounding box detection capabilities to identify query-relevant objects and generate corresponding query-relevant notes. These detections are incorporated into the scene graph by associating them with existing tracks or creating new ones. The 304

362

363

364

365

366

367

368

369

370

query-relevant notes are used to update the correspond-ing scratch-pad entries.

#### 307 analyze\_objects

For each node in the user-supplied list that is visible in 308 309 the frame, the VLM analyzes its appearance and answers query, storing the result in the node's notes. 310 The function employs the VLM to examine each visi-311 ble node's bounding box and generate descriptive notes 312 pertaining to the query. These notes are stored in the 313 corresponding nodes' scratch-pad entries. If specified 314 nodes are not visible in the frame, the function defaults 315 to find\_objects. 316

## 317 analyze\_frame

318A frame-level variant that jointly discovers undetected319objects and annotates existing ones with respect to320query. This consolidated approach can both identify321new perceptual elements and enrich the semantic under-322standing of known objects in a single operation.

These APIs enable a critical capability: the reasoning 323 agent itself can identify and remedy gaps in its scene rep-324 resentation during inference. Rather than preprocessing 325 exhaustively for anticipated questions, GraphPad builds a 326 minimal initial representation and lets task requirements 327 328 guide targeted perceptual refinement. This approach aligns with real-world robotic scenarios where complete scene un-329 derstanding is computationally intractable, but targeted per-330 ception can efficiently support specific goals. 331

By invoking these functions, the agent patches omissions in its memory in real time, yielding a task-conditioned
graph that improves both answer accuracy and confidence
without requiring rescanning of the physical scene.

# **336 4. Results**

We evaluate GraphPad on the OpenEQA benchmark [15] to
assess how language-guided scene graph updates affect spatial reasoning performance. Our experiments use Gemini
2.0 Flash [17] as both the reasoning agent and detector. All
evaluations use the episodic-memory variant of OpenEQA,
where systems receive fixed keyframes from 3D scene scans
(HM3D [19] and ScanNet [4]).

## 344 4.1. Baseline Comparison on OpenEQA

Table 1 presents GraphPad's performance against estab-345 lished baselines. Using an initial frame memory size of 346 347  $n_{img} = 5$ , search depth of m = 20, and Frame-Level API, 348 GraphPad achieves 55.3% accuracy. This represents a 3.0 349 percentage point increase over using the same VLM (Gemini 2.0 Flash) with image-only input (52.3%). GraphPad 350 processes 5 initial frames compared to the 25 frames (every 351 *k*-th frame with k = 5 for HM3D and k = 20 for ScanNet) 352 353 used in the image-only baseline.

Table 1. OpenEQA performance comparison across methods

Method	Accuracy (%)			
Blind LLMs				
GPT-4	33.5			
Socratic LLMs w/ Frame Captions				
GPT-4 w/ LLaVA-1.5	43.6			
Socratic LLMs w/ Scene-Graph Captions				
GPT-4 w/ CG	36.5			
GPT-4 w/ SVM	38.9			
Multi-Frame VLMs				
GPT-4V	55.3			
Gemini-2.0 Flash	52.3			
3D-Mem (w/ GPT4V)	57.2			
Human Agent				
Human	86.8			
Our Results				
GraphPad (w/ Gemini 2.0 Flash)	55.3			
	Method nd LLMs GPT-4 cratic LLMs w/ Frame Captions GPT-4 w/ LLaVA-1.5 cratic LLMs w/ Scene-Graph Cap GPT-4 w/ CG GPT-4 w/ SVM Ilti-Frame VLMs GPT-4V Gemini-2.0 Flash 3D-Mem (w/ GPT4V) Iman Agent Human r Results GraphPad (w/ Gemini 2.0 Flash)			

Table 2. Ablation study of system components

Method	Accuracy (%)
Frame Memory	32.9
Frame Memory + Scene Graph	34.6
Frame Memory + Navigation Log	42.5
Frame Memory + SG + Navigation Log	46.9
Frame Memory + Navigation Log + Image API	45.1
Frame Memory + SG + Navigation Log + Node-level API	47.1
Frame Memory + SG + Navigation Log + Frame-level API	50.5

GraphPad scores higher than static scene graph methods354(CG: 36.5%, SVM: 38.9%), suggesting potential benefits355of dynamic scene graph updates during inference. The system performs identically to GPT-4V (55.3%) despite using357a different base VLM. 3D-Mem [29] achieves 57.2% accuracy, outperforming our approach by 1.9 percentage points,359though it uses GPT-4V rather than Gemini.360

The comparison suggests that structured memory with targeted updates can help bridge the performance gap between different vision-language models while potentially reducing the number of frames that need to be processed.

#### 4.2. Component Analysis

To understand the contribution of different GraphPad components, we conducted an ablation study (Table 2) using a subset of 184 OpenEQA questions. Starting with just four raw frames (32.9% accuracy), we progressively added components to measure their individual impact.

The Navigation Log provides the largest individual improvement (+9.6 percentage points over frame-only), suggesting that structured information about frame contents371and the largest individual improvement (+9.6 percentage points over frame-only), suggesting that structured information about frame contents372and the largest individual improvement (+9.6 percentage points over frame-only), suggesting that structured information about frame contents372and the largest individual improvement (+9.6 percentage points over frame-only), suggesting that structured information about frame contents373and the largest individual improvement (+9.6 percentage points over frame-only), suggesting that structured information about frame contents373and the largest individual improvement (+9.6 percentage points over frame-only), suggesting that structured information about frame contents373and the largest individual improvement (+9.6 percentage points over frame-only), suggesting that structured information about frame contents374

Table 3. Performance as a function of search depth

Search Depth	Accuracy (%)	
0	46.9	
1	45.8	
2	46.3	
3	48.1	
4	45.9	
5	49.7	
20	51.2	

Table 4. Performance as a function of initial frame count

Initial Frames	Accuracy (%)	
2	48.2	
3	47.4	
4	49.9	
5	50.3	
6	50.5	

graph yields a modest improvement (+1.7 percentage points
over frame-only), while the Modifiability APIs add another
3.6 percentage points over the static scene representation.

Frame-level APIs (50.5%) outperformed node-level variants (47.1%). We observed that when using node-level APIs, the VLM tended to rely more on find\_objects
than analyze\_objects, often using the former in ways similar to analyze\_frame.

# **4.3. Effect of Search Depth and Frame Count**

We examined how search depth (m) affects accuracy using 384 the OpenEQA184 subset with  $n_{ima} = 4$  initial frames and 385 386 the Frame-Level API. Table 3 shows that accuracy generally 387 increases with search depth, reaching 51.2% at m = 20. Interestingly, limited search (m = 1 or m = 2) sometimes 388 performs worse than no search (m = 0), possibly because 389 390 preliminary updates without follow-up refinement can in-391 troduce misleading information.

The initial number of frames in Frame Memory also influences performance (Table 4). Accuracy improves as more frames are included, though gains diminish beyond frames. We observed that with fewer initial frames, the model made more API calls on average, potentially compensating for the limited initial context.

# **398 4.4. Performance by Question Category**

Breaking down performance by question category (Table 5) reveals variation in how GraphPad (with k = 5) compares to using Gemini-2.0 with only images (every 5th frame in the scene).

GraphPad shows larger improvements in attribute recog nition (+20.3 percentage points), functional reasoning (+5.7

Table 5. Performance by question category: GraphPad vs. Gemini 2.0 with images only

Category	GraphPad (%)	Gemini (%)
Attribute Recognition	66.8	46.5
<b>Object State Recognition</b>	69.6	66.5
Functional Reasoning	59.2	53.5
World Knowledge	55.9	52.0
Object Recognition	58.4	62.2
Spatial Understanding	47.7	52.4
Object Localization	31.3	34.3



Figure 2. Distribution of API calls per query on OpenEQA

percentage points), and object state recognition (+3.1 per-<br/>centage points). These categories often require detailed<br/>analysis of specific object properties. The image-only base-<br/>line performs better in object recognition (-3.8 percentage<br/>points), spatial understanding (-4.7 percentage points), and<br/>object localization (-3.0 percentage points), which may rely<br/>more on raw visual processing capabilities.405<br/>406<br/>407

# 4.5. API Call Distribution

Analysis of GraphPad's API usage (Fig. 2) shows that 95%413of questions are answered with 5 or fewer API calls, with an<br/>average of 1.9 calls per question. The distribution indicates<br/>that in most cases, the system requires relatively few up-<br/>dates to answer questions, though we observe a long tail of<br/>more complex queries requiring additional refinement steps.413413

The peak at zero calls represents questions that could be<br/>answered using only the initial scene representation without<br/>any additional API calls, suggesting that for some question<br/>types, the initial structured memory provides sufficient con-<br/>text.419<br/>420<br/>421

# 5. Limitations

While GraphPad demonstrates the value of editable 3D425scene representations, several important limitations affect426its current implementation:427

• Error propagation in detection. Our current implementation lacks a verification mechanism for object detection 428

412

424

493

494

495

496

497

498

499

500

501

502

503

504

505

506

507

508

509

510

511

512

513

514

515

516

517

518

519

520

521

522

523

524

525

526

527

528

529

530

531

532

533

534

535

quality. When the VLM misidentifies an object or relation, this error becomes part of the scene graph and can
negatively affect downstream reasoning.

- API design constraints. The three operations (find, analyze objects, analyze frame) represent our initial attempt at a minimal API set, but we have not rigorously evaluated whether this is optimal. Different reasoning tasks might benefit from specialized APIs not explored in this work.
- Computational overhead. Each API call requires a full
  VLM inference pass, adding significant latency (typically
  2-3 seconds per call in our implementation). This latency
  currently limits GraphPad's applicability to real-time systems.
- Domain generalization. Our APIs were designed specifically for question answering about static scenes. We have not tested their effectiveness for other domains like manipulation planning or navigation in dynamic environments.
- Scalability limits. As scene graphs grow larger, both the prompt size and reasoning complexity increase. Our experiments were limited to medium-sized home environments; performance in larger spaces remains untested.

# 453 6. Conclusion

We presented GraphPad, a system that enables VLMs to
update 3D scene graphs during inference through languagecallable functions. On OpenEQA, GraphPad improved accuracy by 3.0 percentage points over an image-only baseline using the same VLM while requiring fewer input
frames.

Our experiments suggest several insights for 3D 460 language-vision systems. First, static structured represen-461 tations appear to benefit from targeted, task-specific re-462 finement. Second, language-guided perception may of-463 fer a middle ground between exhaustive preprocessing and 464 purely reactive vision. Third, the efficiency of GraphPad 465 (averaging under 2 API calls per question) indicates that 466 targeted scene exploration can be a practical strategy. 467

For future 3D vision-language-action systems, these results suggest investigating how reasoning agents might direct their own perception in service of task goals. Extending GraphPad's approach to manipulation planning, navigation, and dynamic scenes could help bridge the gap between
language understanding and effective action in 3D environments.

# 475 **References**

- 476 [1] M. Abate N. Hughes Y. Chang J. Shi A. Gupta L. Carlone
  477 A. Rosinol, A. Violette. Kimera: from SLAM to spatial per478 ception with 3D dynamic scene graphs. In *arxiv*, 2021. 2
- 479 [2] Abrar Anwar, John Welsh, Joydeep Biswas, Soha Pouya, and480 Yan Chang. Remembr: Building and reasoning over long-

horizon spatio-temporal memory for robot navigation. *arXiv* 481 preprint arXiv:2409.13682, 2024. 1, 2 482

- [3] Hao-Tien Lewis Chiang, Zhuo Xu, Zipeng Fu. 483 Mithun George Jacob, Tingnan Zhang, Tsang-Wei Ed-484 ward Lee, Wenhao Yu, Connor Schenck, David Rendleman, 485 Dhruy Shah, Fei Xia, Jasmine Hsu, Jonathan Hoech, Pete 486 Florence, Sean Kirmani, Sumeet Singh, Vikas Sindhwani, 487 Carolina Parada, Chelsea Finn, Peng Xu, Sergey Levine, and 488 Jie Tan. Mobility vla: Multimodal instruction navigation 489 with long-context vlms and topological graphs. arXiv 490 preprint arXiv:2407.07775, 2024. 2 491
- [4] Angela Dai, Angel X. Chang, Manolis Savva, Maciej Halber, Thomas Funkhouser, and Matthias Nießner. Scannet: Richly-annotated 3d reconstructions of indoor scenes. In *Proc. Computer Vision and Pattern Recognition (CVPR), IEEE*, 2017. 5
- [5] Abhishek Das, Samyak Datta, Georgia Gkioxari, Stefan Lee, Devi Parikh, and Dhruv Batra. Embodied question answering. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018. 1
- [6] Runyu Ding, Jihan Yang, Chuhui Xue, Wenqing Zhang, Song Bai, and Xiaojuan Qi. Pla: Language-driven openvocabulary 3d scene understanding. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 7010–7019, 2023. 2
- [7] Elias Greve, Martin Büchner, Niclas Vödisch, Wolfram Burgard, and Abhinav Valada. Collaborative dynamic 3d scene graphs for automated driving. In 2024 IEEE International Conference on Robotics and Automation (ICRA), pages 11118–11124, 2024. 2
- [8] Qiao Gu, Ali Kuwajerwala, Sacha Morin, Krishna Murthy Jatavallabhula, Bipasha Sen, Aditya Agarwal, Corban Rivera, William Paul, Kirsty Ellis, Rama Chellappa, et al. Conceptgraphs: Open-vocabulary 3d scene graphs for perception and planning. In 2024 IEEE International Conference on Robotics and Automation (ICRA), pages 5021–5028. IEEE, 2024. 2
- [9] Daniel Honerkamp, Martin Büchner, Fabien Despinoy, Tim Welschehold, and Abhinav Valada. Language-grounded dynamic scene graphs for interactive object search with mobile manipulation. *IEEE Robotics and Automation Letters*, 2024.
- [10] Chenguang Huang, Oier Mees, Andy Zeng, and Wolfram Burgard. Visual language maps for robot navigation. In Proceedings of the IEEE International Conference on Robotics and Automation (ICRA), London, UK, 2023. 2
- [11] Nathan Hughes, Yun Chang, and Luca Carlone. Hydra: A real-time spatial perception system for 3d scene graph construction and optimization. *arXiv preprint arXiv:2201.13360*, 2022. 1, 2
- [12] Wenhao Li, Zhiyuan Yu, Qijin She, Zhinan Yu, Yuqing Lan, Chenyang Zhu, Ruizhen Hu, and Kai Xu. Llm-enhanced scene graph learning for household rearrangement. In *SIG-GRAPH Asia 2024 Conference Papers*, pages 1–11, 2024. 1, 2
- [13] Peiqi Liu, Zhanqiu Guo, Mohit Warke, Soumith Chintala, Chris Paxton, Nur Muhammad Mahi Shafiullah, and Lerrel
   537

607

608

609

610

611

612

613

614

615

616

617

538Pinto. Dynamem: Online dynamic spatio-semantic mem-539ory for open world mobile manipulation. arXiv preprint540arXiv:2411.04999, 2024. 2

- [14] Dominic Maggio, Yun Chang, Nathan Hughes, Matthew
  Trang, Dan Griffith, Carlyn Dougherty, Eric Cristofalo,
  Lukas Schmid, and Luca Carlone. Clio: Real-time
  task-driven open-set 3d scene graphs. *arXiv preprint arXiv:2404.13696*, 2024. 2
- 546 [15] Arjun Majumdar, Anurag Ajay, Xiaohan Zhang, Pranav
  547 Putta, Sriram Yenamandra, Mikael Henaff, Sneha Silwal,
  548 Paul Mcvay, Oleksandr Maksymets, Sergio Arnaud, et al.
  549 Openeqa: Embodied question answering in the era of foun550 dation models. In *Proceedings of the IEEE/CVF Conference*551 *on Computer Vision and Pattern Recognition*, pages 16488–
  552 16498, 2024. 1, 2, 3, 5
- [16] Songyou Peng, Kyle Genova, Chiyu Jiang, Andrea
  Tagliasacchi, Marc Pollefeys, Thomas Funkhouser, et al.
  Openscene: 3d scene understanding with open vocabularies.
  In Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, pages 815–824, 2023. 2
- [17] Sundar Pichai, Demis Hassabis, and Koray Kavukcuoglu. In troducing Gemini 2.0: our new AI model for the agentic era,
   2024. 5
- [18] Abhinav Rajvanshi, Karan Sikka, Xiao Lin, Bhoram Lee,
  Han pang Chiu, and Alvaro Velasquez. Saynav: Grounding
  large language models for dynamic planning to navigation
  in new environments. In *34th International Conference on Automated Planning and Scheduling*, 2024. 2
- [19] Santhosh Kumar Ramakrishnan, Aaron Gokaslan, Erik Wi-566 567 jmans, Oleksandr Maksymets, Alexander Clegg, John M 568 Turner, Eric Undersander, Wojciech Galuba, Andrew West-569 bury, Angel X Chang, Manolis Savva, Yili Zhao, and Dhruv Batra. Habitat-matterport 3d dataset (HM3d): 1000 large-570 571 scale 3d environments for embodied AI. In Thirty-fifth Con-572 ference on Neural Information Processing Systems Datasets 573 and Benchmarks Track, 2021. 5
- [20] Krishan Rana, Jesse Haviland, Sourav Garg, Jad AbouChakra, Ian Reid, and Niko Suenderhauf. Sayplan: Grounding large language models using 3d scene graphs for scalable
  task planning. In 7th Annual Conference on Robot Learning,
  2023. 2
- 579 [21] Saumya Saxena, Blake Buchanan, Chris Paxton, Bingqing
  580 Chen, Narunas Vaskevicius, Luigi Palmieri, Jonathan Francis, and Oliver Kroemer. Grapheqa: Using 3d semantic scene
  582 graphs for real-time embodied question answering. *arXiv*583 *preprint arXiv:2412.14480*, 2024. 3
- [22] Rutav Shah, Albert Yu, Yifeng Zhu, Yuke Zhu, and Roberto
  Martín-Martín. Bumble: Unifying reasoning and acting with
  vision-language models for building-wide mobile manipulation. arXiv preprint arXiv:2410.06237, 2024. 2
- 588 [23] Ayca Takmaz, Alexandros Delitzas, Robert W Sumner,
  589 Francis Engelmann, Johanna Wald, and Federico Tombari.
  590 Search3d: Hierarchical open-vocabulary 3d segmentation.
  591 arXiv preprint arXiv:2409.18431, 2024. 2
- [24] Yujie Tang, Meiling Wang, Yinan Deng, Zibo Zheng,
  Jingchuan Deng, and Yufeng Yue. Openin: Open-vocabulary
  instance-oriented navigation in dynamic domestic environments. arXiv preprint arXiv:2501.04279, 2025. 2

- [25] Zixuan Wang, Bo Yu, Junzhe Zhao, Wenhao Sun, Sai Hou, Shuai Liang, Xing Hu, Yinhe Han, and Yiming Gan. Karma: Augmenting embodied ai agents with long-and-short term memory systems. arXiv preprint arXiv:2409.14908, 2024.
  2
- [26] Abdelrhman Werby, Chenguang Huang, Martin Büchner, Abhinav Valada, and Wolfram Burgard. Hierarchical openvocabulary 3d scene graphs for language-grounded robot navigation. In *First Workshop on Vision-Language Models for Navigation and Manipulation at ICRA 2024*, 2024. 2, 4
- [27] Shitao Xiao, Zheng Liu, Peitian Zhang, Niklas Muennighoff, Defu Lian, and Jian-Yun Nie. C-pack: Packed resources for general chinese embeddings, 2024. 4
- [28] Quanting Xie, So Yeon Min, Tianyi Zhang, Kedi Xu, Aarav Bajaj, Ruslan Salakhutdinov, Matthew Johnson-Roberson, and Yonatan Bisk. Embodied-rag: General non-parametric embodied memory for retrieval and generation. arXiv preprint arXiv:2409.18313, 2024. 2
- [29] Yuncong Yang, Han Yang, Jiachen Zhou, Peihao Chen, Hongxin Zhang, Yilun Du, and Chuang Gan. 3d-mem: 3d scene memory for embodied exploration and reasoning, 2024. 2, 3, 5
- [30] Mike Zhang, Kaixian Qu, Vaishakh Patil, Cesar Cadena, and Marco Hutter. Tag map: A text-based map for spatial reasoning and navigation with large language models. arXiv preprint arXiv:2409.15451, 2024. 2
  621